GeVLM: 3D Object Grounding with Geometry-Enhanced Vision Language Model

Anonymous ACL submission

Abstract

Understanding 3D scenes with point cloud data in tasks such as object referencing, questionanswering, and captioning poses significant challenges to vision language models (VLMs), due to the complexity of integrating both linguistic and spatial information. While existing methods have mapped point cloud features into LLM space to enable 3D scene comprehension, they often overlook viewpoint information and the relative spatial distance between objects, this can lead to confusion in interpreting spatial descriptions and grounding objects. This paper presents a geometry-enhanced vision LM (GeVLM) to address these challenges. Specifically, we propose viewpoint-consistent position encoding (VCPE) to enhance the relative spatial relationship representation agnostic to the camera viewpoint, and propose the distanceaware cross-entropy (DACE) loss to incorporate distance information in the label space. We additionally introduce the DetailedScanRefer dataset, which provides identifiers and spatial annotation for each object mentioned in the referencing description to further emphasize spatial relationships. GeVLM demonstrates significant improvements over the strong Chat-scene baseline, particularly with 1.3% Acc@0.25 and 1.0% Acc@0.50 improvements on the multiple object setup and state-of-the-art overall performance on ScanRefer dataset¹.

1 Introduction

004

011

017

019

021

034

039

The rapid advancement of Multimodal Large Language Models (LLMs) has greatly enhanced their capabilities in addressing a wide range of tasks involving complex input modalities, such as audio (Tang et al., 2024a; Chu et al., 2023; Gong et al., 2024), images (Liu et al., 2024c,b; Li et al., 2023; Bai et al., 2023; Lin et al., 2023; Chen et al., 2023c) and videos (Zhang et al., 2023a; Cheng et al., 2024; Sun et al., 2024). Recent studies have focused on extending the application of LLMs to the understanding of realistic 3D scenes represented by point clouds(Han et al., 2023; Hong et al., 2023; Wang et al., 2023b; Huang et al., 2023a; Chen et al., 2024b,a), enabling these models to perform tasks such as question-answering, object referencing, and captioning for real-world 3D scenes. Specifically, the task of 3D referencing (Chen et al., 2020) requires LLMs to comprehend detailed object descriptions while simultaneously understanding complex 3D scenes to accurately identify the object being referenced. This task presents significant challenges as it requires the understanding of both linguistic and spatial information. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

081

Previous work in this area has successfully grounded LLMs on 3D point clouds, demonstrating scene comprehension abilities (Hong et al., 2023; Han et al., 2023). Chat-3D (Wang et al., 2023b) maps 3D features into the LLM space and uses a relation module to capture spatial relationships, showcasing strong conversational abilities within 3D environments. Further advancements (Huang et al., 2023a) enhance 3D object referencing by integrating unique identifiers with detailed scene annotations. Nevertheless, these approaches overlook the importance of viewpoint consistency across different examples by simply using 3D world coordinates as input, despite the relative ease of recovering viewpoint information with modern techniques such as SLAM (Thrun, 2008). Moreover, the training objectives are simply cross-entropy (CE) which penalizes other objects equally regardless of whether they are close or far to the target. These factors limit the model performance in scene understanding and object grounding.

To address the aforementioned deficiencies, this paper proposes a geometry-enhanced visual language model (GeVLM) to improve 3D object grounding performance from perspectives including model structure, training criteria and dataset an-

¹We have made all the code, model checkpoints, and DetailedScanrefer used in this work available at https://anonymous.4open.science/r/GeVLM-1372/

notations. Specifically, we propose the viewpointconsistent position encoding (VCPE) which allows 083 relative spatial relationships, e.g. left/right, to be correctly referred to under arbitrary camera viewpoint. Besides, we propose distance-aware crossentropy (DACE) loss which incorporates relative 087 distance information into the label space so that non-target object tokens receive different levels of penalization depending on their spatial affinity to the target. To further boost 3D grounding, we propose the DetailedScanRefer dataset which includes the object identifier and the location for each object mentioned in the description. As a result, GeVLM showed consistent improvements over the strong Chat-scene baseline on a range of 3D scene understanding tasks. The main contributions of this paper are summarized as follows.

> • This paper proposes GeVLM, a geometryenhanced VLM for 3D object referencing and understanding, leveraging easy-to-refer object identifiers. To the best of our knowledge, GeVLM is the first visual LLM that formally investigates and incorporates 3D viewpoint information and relative 3D spatial distance in visual LLMs.

• We propose VCPE to ensure viewpoint consistency in position encoding of point cloud coordinates. In addition, we propose DACE to inject distance information into label space for improved grounding. We also curate the Detailed-ScanRefer dataset with fine-grained identifier annotations for each object in the description.

GeVLM demonstrated state-of-the-art object referencing performance with consistent improvements over the Chat-scene baseline across various 3D scene understanding tasks. In particular, GeVLM achieved 1.3% and 1.0% absolute improvements in Acc@0.25 and Acc@0.50 on the multi-object partition of the ScanRefer dataset.

2 Related Work

100

101

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129

3D Grounding using Language Models Recent research has explored the integration of Large Language Models (LLMs) with 3D object understanding for various applications. LLM-Grounder (Yang et al., 2024) utilizes LLMs to decompose complex queries and evaluate spatial relations for zero-shot 3D visual grounding. Grounded 3D-LLM (Chen et al., 2024b) introduces scene referent tokens and contrastive language-scene pre-training to unify various 3D vision tasks within a generative framework. Transcrib3D (Fang et al., 2024) brings together 3D detection methods and the emergent reasoning capabilities of large language models (LLMs). Cube-LLM (Cho et al., 2024), a multimodal large language model, can ground and reason about 3D objects in images without 3D-specific architectural design or training.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

Language-Driven 3D Scene Understanding There has been growing interest in using natural language to enhance how computers interpret and interact with 3D environments. It involves training models to understand 3D scenes based on verbal instructions. Several tasks were developed to test the understanding ability. Specifically, 3D Visual Grounding (Chen et al., 2020; Huang et al., 2022; Wang et al., 2023a; Hsu et al., 2023; Yang et al., 2024; Unal et al., 2024) involves models identifying a specified object within a 3D scene according to a language query.

3D Large Multi-modal Models Through the usage of large scale 3D object datasets (Yu et al., 2022; Xue et al., 2023; Zhou et al., 2023), 3D Object-level Large Multi-modal Models (LMMs)(Xu et al., 2023; Liu et al., 2024a; Qi et al., 2024; Tang et al., 2024b) have managed to bridge the gap between 3D modality and texts. However, these models fall short when complex spatial reasoning is needed for 3D scenes. Therefore, multiple models (Ziyu et al., 2023; Wang et al., 2023b; Huang et al., 2023a; Chen et al., 2024b) have been proposed as scene-level LLMs. 3D-LLM (Hong et al., 2023) uses point clouds and their features as input and can handle various 3D-related tasks. The model attempts to improve the understanding of complex spatial relationships among objects by using positional embeddings and learning location tokens. However, the model projects 3D features into the 2D feature space of a pretrained visionlanguage model, posing significant challenges to capture the 3D spatial structure and complex relationships among objects. Chat-3D (Wang et al., 2023b) and Chat-Scene (Huang et al., 2023b) directly utilizes 3D scene-text data to align the 3D scene with large language model (Llama). However, Chat-3D could only handle one target object per conversation. To overcome this limitation, Chat-3D-v2 (Huang et al., 2023a), as our baseline model, introduced unique object identifiers in addition to 3D object features, and significantly improved the 3D grounding performance.

247

248

249

250

251

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

231

232

3 Geometry-Enhanced Visual LM

181

182

183

188

189

191

193

195

196

197

198

199

203

204

207

208

210

We introduce a novel geometry-enhanced visionlanguage model (GeVLM) specifically designed for 3D object grounding tasks. As illustrated in Figure 1(A), GeVLM integrates several components: a 3D segmenter, a 3D feature encoder, a 2D multiview image encoder, a 3D position encoder, and a pretrained LLM. In this setup, the 3D segmenter and both the 3D and 2D feature encoders remain frozen during training. Our primary objective is to fine-tune the pretrained LLM so it can effectively interpret language-based referring expressions by incorporating 3D geometric cues.

To validate the versatility of our approach, we build GeVLM on top of two baseline frameworks: Chat-3D-v2 and Chat-scene. The key difference between the two frameworks is that Chat-scene includes multi-view object images and employs a 2D encoder, while Chat-3D-v2 does not.

These geometric cues are considered in two key aspects. First, we propose a viewpoint-consistent position encoding (VCPE) to account for camera perspective in 3D scene understanding, as detailed in Section 3.1. Second, we introduce a distanceaware cross-entropy (DACE) loss, discussed in Section 3.2, to highlight the importance of spatial affinity in the grounding task. Additionally, in Section 4, we present a densely annotated grounding dataset, curated with assistance from GPT-40.

3.1 Viewpoint-Consistent Position Encoding

Imagine you are inside a room and refer to the 211 chair in front of you. The success of the refer-212 213 ring depends on the viewpoint of the observer. In other words, ambiguities will arise if the viewpoint 214 is not known. This is supported by the fact that 215 the annotation from ScanRefer includes the cam-216 era pose information. However, existing methods 217 (Wang et al., 2023b; Huang et al., 2023a) overlook the viewpoint information, hence refer to the same 219 3D point cloud when querying different referring descriptions. We notice that the incorrect ground-221 ing outcomes are mainly due to the rotation of the camera viewpoint, which makes relative spatial descriptions such as left/right and front/back confusing to LLMs. For example, in the 3D scenein Fig.2 with 4 different viewpoints, the description "the shelf is to the right of the bed" only makes sense when observing the scene from a consistent viewpoint, e.g. 1 and 4. Nevertheless, methods like Chat-3D and 3D-LLM ignore camera viewpoint, 230

and directly utilize world coordinates as input for object grounding. This inevitably introduces viewpoint inconsistency to the model training and leads to sub-optimal performance.

In GeVLM, we carefully transform 3D point cloud to ensure viewpoint consistency across referring expressions. Based on the transformed coordinates, we propose a position encoding module, VCPE, to effectively learn the relative spatial relationship for downstream 3D tasks. Specifically, to achieve viewpoint consistency, we apply a 3D transformation using the rotation matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ from the camera's extrinsic parameters. The translation vector is omitted to maintain a consistent scene scale across different datasets and tasks. For an object, its centre point $\mathbf{v} \in \mathbb{R}^3$ is transformed to $\mathbf{v}_{rot} = \mathbf{R}\mathbf{v}$. This transformation preserves the spatial configuration of objects relative to the camera orientation and aligns with the viewpoint-dependent language description. As a result, VCPE is crucial for VLMs to effectively generalize across varying viewpoints.

To capture complex spatial relationships, we apply Fourier Feature Mapping (Tancik et al., 2020) to map the low-dimensional coordinates v_{rot} to capture high-frequency details as shown in Eqn (1):

 $\gamma(\mathbf{v}_{\text{rot}}) = \left[\sin(2\pi \mathbf{B} \mathbf{v}_{\text{rot}}), \cos(2\pi \mathbf{B} \mathbf{v}_{\text{rot}})\right], \quad (1)$

where $\mathbf{B} \in \mathbb{R}^{3 \times D}$ is a Gaussian random matrix, and D is the dimensionality of the Fourier features. This mapping projects the rotated points into a higher-dimensional space, enabling the model to represent positional information with high-frequency components. The Fourier embeddings for all objects are concatenated into a matrix $\mathbf{F}^{F} \in \mathbb{R}^{O \times 2D}$, where O represents the number of objects. These embeddings are then projected through a linear layer, followed by a Gaussian Error Linear Unit (GELU) activation and a multi-head self-attention layer, as shown in Eqn (2):

$$\mathbf{F}^{\text{attn}} = \text{MHSA}(\text{GELU}(\mathbf{F}^{\text{F}}\mathbf{W} + \mathbf{b})),$$
 (2)

where $\mathbf{W} \in \mathbb{R}^{2D \times D'}$ and $\mathbf{b} \in \mathbb{R}^{D'}$ are learnable parameters, MHSA(·) denotes multi-head selfattention, and D' is the dimensionality of the projected features. This produces attention-weighted embeddings $\mathbf{F}^{\text{attn}} \in \mathbb{R}^{O \times D'}$ that effectively capture spatial relationships. By integrating these components, VCPE improves the model's capacity to comprehend complex spatial configurations.



Figure 1: The model structure of the GeVLM (A) together with distance-aware CE loss (B) and an illustration of the viewpoint consistent position encoding (C).



Description: It's a black three levels shelf. It is located to the right of the bed.

Figure 2: Example scene where viewpoint consistency is vital. The target shelf in the description is only **right** to the bed in the first viewpoint, and the description confuses the model when using other viewpoints, resulting in an incorrect grounding outcome.

3.2 Distance-Aware Cross-Entropy Loss

279

281

282

294

298

Most 3D VLMs commonly rely on language loss to fine-tune 3D tasks due to its simplicity. Efforts have been made to unify multimodal tasks under a single language-based objective. However, we argue that applying standard cross-entropy (CE) loss to 3D grounding tasks is inadequate. Specifically, when training a model with CE to predict the token for a referred object, it penalizes all other object tokens equally. This contrasts with 3D object detection and segmentation, where the goal is to minimize the distance between the ground truth and predictions.

Building on this insight, we propose DACE to incorporate geometric distance between objects into the loss computation. This approach allows spatial relationships to be considered during training. We categorize tokens into *regular tokens* and *object tokens*. We append 100 object tokens to represent scene objects to existing tokens. For instance, object token <OBJ000> will be indexed by 32,000 in Llama2. The DACE loss differentiates regular tokens and object token: standard CE loss is applied to the regular tokens, while a soft label is used for the object token predictions. We further detail the DACE loss next.

299

300

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

322

As shown in part (B) of Fig.1, for each scene, we precompute a distance matrix $\mathbf{D} \in \mathbb{R}^{\mathcal{V}_{obj} \times \mathcal{V}_{obj}}$ where D_{ij} denotes the Euclidean distance between objects *i* and *j*, and \mathcal{V}_{obj} is the number of object tokens. Then, the DACE loss is defined in Eqn. (3).

$$\mathcal{L}_{\text{dist}} = \frac{1}{L} \sum_{i=1}^{L} \boldsymbol{m}_i \cdot \text{CE}(\mathbf{w}_i, P_{\theta}(\mathbf{y}_i | X_i)) + (1 - \boldsymbol{m}_i) \cdot \text{CE}(\hat{\mathbf{y}}_i, P_{\theta}(\mathbf{y}_i | X_i)),$$
(3)

where $m_i = 1$ for object tokens and $m_i = 0$ for regular tokens. CE(·) denotes the cross-entropy loss, L is the total number of tokens in the sequence, and \hat{y}_i is the one-hot label vector. The distanceaware soft label, w_i , is computed in Eqn. (4).

$$\boldsymbol{w}_{i} = \exp(-D_{ij}/T) / \sum_{k=1}^{\mathcal{V}_{obj}} \exp(-D_{ik}/T) \quad (4)$$

where $D_{ij} \in [0, 1]$ is the min-max normalized distance between object *i* and object *j*, and *T* is the temperature parameter controlling the sharpness of the soft label. The intuition behind this loss is to encourage the model to focus on objects with close affinity, rather than on more distant yet semantically similar objects. This is particularly useful

358



Figure 3: Left: DACE selects the correct chair among the confusing row of chairs based on the location information. Right: DACE accurately finds the sofa near the coffee table rather than other similar sofas.

in scenarios with multiple similar objects, such as chairs in a meeting room, where a specific chair is being referred to, as shown in Fig.3.

4 DetailedScanRefer: A Densely Annotated Grounding Dataset

323

324

325

327

329

334

335

337

339

340

343

345

347

352

354

357

To further improve grounding, we introduce the DetailedScanRefer dataset, an extension of the Scan-Refer dataset (Chen et al., 2020). DetailedScanRefer features annotations for both target and landmark objects, in contrast to Chat-3D-v2 (Huang et al., 2023a) which only annotates the target. Unique object identifiers (object IDs) are assigned to each object in the 3D scene. The generation pipeline is shown in Fig.4. All objects in the description are matched to IDs from Mask3D (Schult et al., 2023) for consistent object referencing.

Scene Image Retrieval via Camera Pose Matching Due to the poor quality of images directly rendered from the 3D point, we retrieve a photo of the real-world scene with the most similar view from the ScanNet dataset (Dai et al., 2017)². For each description in ScanRefer, we retrieve the closest camera pose from ScanNet, along with its corresponding RGB image and depth map. The best matching camera pose $T_{best} \in \mathbb{R}^{3\times 4}$ is determined by minimizing the mean Euclidean distance between the camera coordinates of the entire scene:

$$\mathbf{T}_{\mathsf{best}} = rg\min_{\mathbf{T}_i} \left(rac{1}{N} \sum_{k=1}^N \left\| oldsymbol{p}_k(\mathbf{T}_{\mathsf{target}}) - oldsymbol{p}_k(\mathbf{T}_i)
ight\|_2
ight)$$

where N is the total number of points in the scene, $\mathbf{T}_{target} \in \mathbb{R}^{3 \times 4}$ is the camera pose corresponding to the ScanRefer description, and $\mathbf{T}_i \in \mathbb{R}^{3 \times 4}$ is the *i*-th candidate pose from the same scene in ScanNet. The term $\mathbf{p}_k(\mathbf{T}) \in \mathbb{R}^3$ represents the camera coordinates of the k-th point in the scene transformed by the camera pose **T**. The Euclidean distance between the transformed points under different camera poses is averaged over all points in the scene, ensuring that the selected camera pose closely matches the viewpoint described by the ScanRefer description.

Visibility and Object Annotation Simply projecting an object's center onto the image can lead to incorrect annotations, as hidden or partially visible objects may be included by mistake. This becomes particularly problematic in later stages, such as querying GPT40 for high-quality responses, where accurately labeling only visible objects is crucial. To address this issue, we project the 3D instance segmentation mask to image space and compare it with the scene's depth map. We project 3D points onto the 2D image plane using the camera parameters as follows:

$$u = \frac{f_x X_c}{Z_c} + c_x, \quad v = \frac{f_y Y_c}{Z_c} + c_y$$
37

where (X_c, Y_c, Z_c) represent the 3D point in the camera coordinate system, and f_x, f_y and c_x, c_y denote the focal lengths and principal points of the camera, respectively. Visibility is confirmed by comparing the estimated depth Z_c with the depth map D_{depth} , using the condition visible = $|Z_c - D_{depth}(u, v)| \leq \delta$, where δ =0.1 meter accounts for minor discrepancies due to sensor noise. Appendix A.2 shows the pixel-level visibility masks and how objects are annotated.

Photo Annotation We generate annotated images by overlaying unique object identifiers at the mean pixel coordinates of each object's mask. For example, an object with index 13 is labelled as "OBJ013" to clearly tag visible objects in the image. Examples of these annotations can be found in the bottom row of Fig.6. These annotated images are then sent to the GPT-4 API, along with the original ScanRefer description, for automatic generation of detailed annotations. In the generated descriptions, object IDs are inserted after the object references. As shown in Fig.4, for the original description: "This is a brown guitar. It is leaning against the wall." The enhanced output is: "This is a brown guitar **<OBJ018>**. It is leaning against the wall <OBJ032>." Details of the prompts can be found in Appendix A.3.

Data Cleaning and Quality Rating To ensure high-quality annotations, we implemented several data cleaning processes. Key steps included discarding annotations where the first object ID did

²Examples are shown in Appendix A.1



Figure 4: DetailedScanRefer generation pipeline. Given a ScanRefer description, we first retrieve its corresponding camera pose, T_{target} . Using a camera pose matching algorithm, we find the closest match, T_{best} , from the ScanNet dataset, along with the corresponding image I_{best} and depth map D_{best} . The semantic segmentation result S is then projected from 3D space onto the image using T_{best} and the intrinsic matrix of the scene. D_{best} is applied to filter visible pixels for each object, and the visible object IDs are annotated on I_{best} . Finally, GPT-40 is used to append an object ID to each object in the description.

not match the ground truth, in line with ScanRefer's assumption that the target object is described first. Additionally, any outputs containing NaN values were removed. The cleaned annotations were then used as ground truth for training, where the model predicted object IDs for each mentioned object. Detailed statistics for each data cleaning step are provided in Appendix A.4. Furthermore, annotation quality evaluation details using GPT-40 are presented in Appendix A.5.

5 Experiments

407

408

409

410 411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

5.1 Experimental Setup

Training Data We follow exactly the same training data setup as Chat3D-v2 (Wang et al., 2023b) so that our results are directly comparable. The training datasets include ScanRefer (Chen et al., 2020), Scan2Cap (Chen et al., 2021), ScanQA (Azuma et al., 2022), SQA3D (Ma et al., 2023), Multi3DRef (Zhang et al., 2023b), and NR3D (Achlioptas et al., 2020). We also use ObjAlign, which is a dataset for aligning object IDs with objects³. Among these datasets, only ScanRefer and Scan2Cap tasks use viewpoint information, as they are the only datasets providing it. The proposed DetailedScanRefer, with about 16,000 samples in total, is also used where specified. For validation, we use ScanRefer, Scan2Cap, ScanQA, SQA3D, and Multi3DRef to select the best model checkpoint. The Scan2Cap dataset is modified by associating a camera pose with each caption under the GeVLM (Chat-3D v2) setup. We refer to Appendix B for details.

Model and Training Specifications To extract object features, we utilize the pretrained Uni3D (Chen et al., 2023a) as the 3D encoder, which is frozen during training. Mask3D is employed for consistent and accurate segmentation of the 3D data and is frozen during training. OpenIns3D (Huang et al., 2024b) was used to assist with visibility checks and to develop visualization tools. Two versions of GeVLM are developed following the Chat-3D-v2 architecture with Vicuna-v1.5-7B backbone and a Chat-scene(Huang et al., 2023b) architecture with Llama-3-Instruct-8B (Meta AI, 2023) backbone respectively. The VCPE module uses a 256-dim final positional embedding and 128dim Fourier features. The positional embedding projection layer and the multi-head attention module are trainable components. There are 100 object proposals for each scene. The entire training process, using the Adam optimizer and a cosine learning rate scheduler for 3 epochs, requires approximately 11 hours on 4 NVIDIA A100 GPUs.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

Evaluation Metrics For 3D grounding tasks, grounding accuracy is measured at two Intersections over Union (IoU) thresholds: 25% and 50%, referred to as Acc@0.25 and Acc@0.50, respectively. For language tasks, metrics such as BLEU scores (Papineni et al., 2002), METEOR (Lavie and Agarwal, 2007) and CIDEr (Vedantam et al., 2015) are used to measure the degree of overlap between the generated answer and the reference, with higher scores indicating better performance.

5.2 Experimental Results

Results on 3D Grounding Tasks We first show the 3D grounding performance using ScanRefer

³Questions are in the format of "what is Obj14" and the answer is "chair".

Table 1: Accuracy on ScanRefer (Chen et al., 2020) validation set using GeVLM at 0.25 and 0.50 IoU. Unique subset contains samples where the grounding object is unique in the scene, in contrast to Multiple where there are multiple objects of the same kind as the grounding object.

Second acres	Uni	que	Mul	tiple	Overall		
System	Acc@0.25	Acc@0.50	Acc@0.25	Acc@0.50	Acc@0.25	Acc@0.50	
3DJCG (Cai et al., 2022)	-	64.3	-	30.8	-	37.3	
D3Net (Chen et al., 2022a)	-	72.0	-	30.1	-	37.9	
3D-LLM (BLIP2-flant5) (Hong et al., 2023)	-	-	-	-	30.3	-	
Chat-3D v2* (Huang et al., 2023a)	79.0	74.5	34.7	31.6	42.9	39.6	
Chat-scene (Huang et al., 2023b)	89.6	82.5	47.8	42.9	55.5	50.2	
GeVLM (Chat-3D v2) (Ours)	82.0	75.7	39.0	34.7	46.9	42.3	
GeVLM (Chat-scene) (Ours)	87.8	80.0	49.1	43.9	56.3	50.6	

Table 2: Performance comparison. "Expert models" are for specific tasks using task-oriented heads, while "Unifies models" are designed for general instructions and responses for multiple tasks.

Catalan	S	ScanF	Refer	Multi3I	ORefer	Scar	n2Cap	ScanQA		SQA3D	
Category	System	Acc@0.25	Acc@0.5	F1@0.25	F1@0.5	C@0.5	B-4@0.5	С	B-4	EM	EM-R
	ScanRefer (Chen et al., 2020)	37.3	24.3	-	-	-	-	-	-	-	-
	ScanQA (Azuma et al., 2022)	-	-	-	-	-	-	64.9	10.1	-	-
	3DJCG (Cai et al., 2022)	49.6	37.3	-	-	49.5	31.0	-	-	-	-
Ermont Modele	3D-VLP (Jin et al., 2023)	51.4	39.5	-	-	54.9	32.3	67.0	11.1	-	-
Expert Models	M3DRef-CLIP (Zhang et al., 2023b)	51.9	44.7	42.8	38.4	-	-	-	-	-	-
	3D-VisTA (Ziyu et al., 2023)	50.6	45.5	-	-	66.9	34.0	72.9	13.1	48.5	-
	ConcreteNet (Unal et al., 2024)	50.6	46.5	-	-	-	-	-	-	-	-
	Vote2Cap-DETR++ (Chen et al.,					67.6	27.1				
	2023b)	-	-	-	-	07.0	37.1	-	-	-	
	3D-LLM(BLIP2-flant5) (Hong et al.,	30.3			_	_	_	69.4	12.0	_	
	2023)	50.5						07.4	12.0		
	LL3DA (Chen et al., 2024a)	-	-	-	-	65.2	36.8	76.8	13.5	-	-
	LEO (Huang et al., 2024a)	-	-	-	-	72.4	38.2	101.4	13.2	50.0	52.4
Unified Models	Scene-LLM (Fu et al., 2024)	-	-	-	-	-	-	80.0	11.7	53.6	-
Unified Widdels	Chat-3D v2 (Huang et al., 2023a)	42.5	38.4	45.1	41.6	63.9	31.8	87.6	14.0	54.7	-
	Chat-scene (Huang et al., 2023b)	55.5	50.2	57.1	52.4	77.1	36.3	87.7	14.3	54.6	57.5
	GeVLM (Chat-3D v2) (Ours)	46.9	42.3	50.0	46.1	-	-	90.5	15.4	53.5	56.0
	GeVLM (Chat-scene) (Ours)	56.3	50.6	57.1	52.4	81.9	38.1	89.7	14.3	56.5	59.4

and Multi3DRefer datasets. On the ScanRefer 472 dataset shown in Tab.1, the proposed GeVLM achieved consistent performance improvement compared to the Chat-3D v2 baseline with clear margins for both Unique and Multiple subsets. The improvement is particularly pronounced when there are multiple confusing objects with similar semantic classes in the scene, demonstrating the importance of viewpoint and relative distance information which are crucial to distinguishing those objects. Overall, GeVLM (Chat-3D v2) achieved a 4.0% absolute accuracy improvement at 0.25 IoU and a 2.7% improvement at 0.50 IoU respectively compared to the Chat-3D v2 baseline.

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

We then extend our experiments to the Multi3DRefer dataset as shown in Tab.8 where an overall 1.2% absolute F1 score improvement is achieved. In particular, large improvements are found when there are semantically distracting classes, with a 5.8% absolute F1 score improvement on the zero target subset (i.e. the target object is not in the scene) and 3.5% on the single target subset when distractors are added. For the MT subset, where multiple objects sharing the same semantics need to be grounded, we observe a 5.4% performance drop compared to the baseline method. This can be attributed to the nature of the task. First, grounding multiple objects requires less spatial reasoning, making our proposed VCPE less effective. Second, since the task involves grounding multiple objects that share the same semantics, the model relies more heavily on object category recognition than on spatial differentiation, further diminishing the effectiveness of the DACE loss.

495

496

497

498

499

500

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

Results on Language Tasks: 3D QA and Captioning In addition to object grounding, GeVLM is also beneficial in language tasks, as shown in Table 9. While not explicitly designed to enhance language tasks, GeVLM achieves the best performance across most metrics compared to other 3D LLMs capable of 3D grounding tasks on ScanQA and achieved on-par performance with Chat-3D v2. Notably, among the listed models, only 3D-LLM, Chat-3D v2, and GeVLM possess grounding capabilities, emphasizing our model's superior ver-

Table 3: Ablation study on ScanRefer validation set for GeVLM (Chat-3D v2) with VCPE, DetailedScanrefer (Detailed), and DACE. "World", "Camera", and "Rotate" refer to world coordinates, camera coordinates (both rotation and translation), and rotated coordinates (rotation only), respectively.

Mathad	VCDE	Detailed	DACE	Unique		Mul	tiple	Overall		
Wiethod	VCFE	Detalleu	DACE	@0.25	@0.50	@0.25	@0.50	@0.25	@0.50	
Chat-3D v2*	-	-	-	79.0	74.5	34.7	31.6	42.9	39.6	
	World	-	_	81.1	76.1	35.8	32.0	44.2	40.2	
	Camera	-	_	79.0	74.1	35.6	32.2	43.6	40.0	
Ouro	Rotate	-	_	79.6	74.7	36.2	32.6	44.2	40.4	
Ours	Rotate	\checkmark	_	80.7	74.9	35.7	32.3	44.0	40.2	
	Rotate	-	T = 0.05	79.5	73.7	37.9	33.7	45.6	41.1	
	Rotate	1	T = 0.05	80.4	74.1	38.1	34.0	46.0	41.4	
	Rotate	\checkmark	<i>T</i> =0.03	82.0	75.7	39.0	34.7	46.9	42.3	

Table 4: Comparison of predicted vs. target object center distances on Scanrefer, Both average and median distances are reported for unique and multiple scenarios, with and without DACE loss. Med. stands for median.

System	Uni	que	Multiple			
-J	Mean	Med.	Mean	Med.		
VCPE(r) VCPE(r) + DACE	0.60 0.52	0.04 0.04	1.61 1.49	1.13 0.91		

satility and performance in both grounding and language-based tasks. Detailed Scan2Cap results are provided in Appendix B.

517

518

519

520

521

522

523

524

526

530

532

534

535

536

538

541

542

Ablation Studies Ablation studies are performed using the ScanRefer dataset to better understand the effect of each proposed component in GeVLM, as shown in Tab.3. The world, camera, and rotated coordinate systems were analyzed using the same VCPE model. The rotated coordinates showed the best performance. VCPE is especially useful in the multi-object case, where it helps clarify which object is being referenced among similar objects by focusing on their relative positions. Accuracy improves further in the Multiple case when using the DACE loss, which emphasizes spatial distance in the label space, rather than relying only on semantic similarity. In the Unique case, using DetailedScanrefer alone improves performance, but the best results are achieved when combining DetailedScanrefer with DACE loss and VCPE.

Moreover, to quantitatively demonstrate the effect of the DACE loss, the average (mean and median) distance between the predicted and the target object centers across all test samples is presented in Tab.4. The significantly smaller average distances observed in both scenarios indicate that the DACE loss helps GeVLM focus on locations that are spatially closer to the target object.

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

Qualitative Analysis We qualitatively demonstrate the advantages of GeVLM through examples in Fig. 8 of Appendix E, highlighting four common types of spatial confusion: left/right, near/far with respect to the camera, front/back, and geographical directions. In cases (A) to (C), GeVLM, equipped with VCPE and DACE, accurately selects the correct object based on the description, effectively resolving ambiguities that the baseline model fails to address. The baseline consistently picks objects with the correct semantic category but incorrect spatial positioning. In case (D), although GeVLM's prediction does not perfectly overlap with the ground truth bounding box, it aligns with the described location, indicating a better understanding of the spatial context compared to the baseline, which selects an object without considering positional cues.

6 Conclusion

This paper introduces a GeVLM to improve 3D object grounding and scene understanding. By integrating the VCPE module and DACE loss, GeVLM achieves improved interpretation of spatial relationships while effectively incorporating distance information into the label space. Additionally, the DetailedScanRefer dataset with dense object annotation is proposed to enhance the model's spatial reasoning capabilities. GeVLM achieves significant performance gains over the strong Chat-scene baseline, particularly with **1.3%** in Acc@0.25 and **1.0%** in Acc@0.50 improvements on the multiobject partition of the ScanRefer benchmark.

7 Limitations

577

604

606

607

610

611

612

613

614

615

616

617

618

619

621

622

623

627

Our work has several limitations, including the need for further refinement of the camera pose-579 matching algorithm to achieve optimal performance. Our approach is specifically tailored to tasks that rely heavily on viewpoint information, 583 making accurate camera pose information essential to fully leverage the model's capabilities. Many existing datasets lack this type of annotation, limiting the applicability of our method. Furthermore, the prompt design for photo annotation could be 587 improved to enhance both efficiency and precision. Lastly, while this paper uses the Chat-3D-v2 as the 589 baseline and follows the exact model and training configuration for direct comparability, we also notice that, as a fast-evolving field, the latest work, 592 such as Chat-Scene (Huang et al., 2023b), has pro-593 posed foundation models that surpass the performance of Chat-3D-v2 baseline by clear margins. However, our proposed methods are orthogonal to these advancements and, in theory, could be applied to achieve further improvements. Exploring these opportunities will be important future work when additional resources and time are available. 600

References

- Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. 2020.
 ReferIt3d: Neural listeners for fine-grained 3d object identification in real-world scenes.
- Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. 2022. ScanQA: 3d question answering for spatial scene understanding. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 19107–19117. IEEE.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-VL: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Daigang Cai, Lichen Zhao, Jing Zhang, Lu Sheng, and Dong Xu. 2022. 3djcg: A unified framework for joint dense captioning and visual grounding on 3d point clouds. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 16443–16452.
- Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. 2020. Scanrefer: 3d object localization in rgb-d scans using natural language. In *Computer Vision–ECCV 2020: 16th European Conference, Glas-*

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

- Dave Zhenyu Chen, Ali Gholami, Matthias Niesner, and Angel X. Chang. 2021. Scan2cap: Contextaware dense captioning in RGB-d scans. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 3192–3202. IEEE.
- Dave Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias Nießner, and Angel X. Chang. 2023a. UniT3d: A unified transformer for 3d dense captioning and visual grounding. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pages 18063–18073. IEEE.
- Dave Zhenyu Chen, Qirui Wu, Matthias Nießner, and Angel X. Chang. 2022a. D3net: A unified speakerlistener architecture for 3d dense captioning and visual grounding. *Preprint*, arXiv:2112.01551.
- Jiaming Chen, Weixin Luo, Xiaolin Wei, Lin Ma, and Wei Zhang. 2022b. Ham: Hierarchical attention model with high performance for 3d visual grounding.
- Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. 2022c. Language conditioned spatial relation reasoning for 3d object grounding. *Preprint*, arXiv:2211.09646.
- Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. 2024a. Ll3da: Visual interactive instruction tuning for omni-3d understanding, reasoning, and planning. In *Proc. CVPR*.
- Sijin Chen, Hongyuan Zhu, Mingsheng Li, Xin Chen, Peng Guo, Yinjie Lei, Gang Yu, Taihao Li, and Tao Chen. 2023b. Vote2cap-detr++: Decoupling localization and describing for end-to-end 3d dense captioning. *Preprint*, arXiv:2309.02999.
- Yilun Chen, Shuai Yang, Haifeng Huang, Tai Wang, Ruiyuan Lyu, Runsen Xu, Dahua Lin, and Jiangmiao Pang. 2024b. Grounded 3d-llm with referent tokens. *arXiv preprint arXiv:2405.10370*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023c. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. 2024. VideoLLaMA 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*.
- Jang Hyun Cho, Boris Ivanovic, Yulong Cao, Edward Schmerling, Yue Wang, Xinshuo Weng, Boyi Li, Yurong You, Philipp Krähenbühl, Yan Wang, et al.

791

792

793

- 695 701 706 710 711 712 714 715 716 717 718 719 724 727 730 731 733 734 735 738

683

684

737

2024. Language-image models with 3d understanding. arXiv preprint arXiv:2405.03685.

- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-Audio: Advancing universal audio understanding via unified large-scale audiolanguage models. arXiv preprint arXiv:2311.07919.
- Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 5828-5839.
- Jiading Fang, Xiangshan Tan, Shengjie Lin, Igor Vasiljevic, Vitor Guizilini, Hongyuan Mei, Rares Ambrus, Gregory Shakhnarovich, and Matthew R Walter. 2024. Transcrib3d: 3d referring expression resolution through large language models. arXiv preprint arXiv:2404.19221.
- Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. 2024. Scene-llm: Extending language model for 3d visual understanding and reasoning. Preprint, arXiv:2403.11401.
- Yuan Gong, Hongyin Luo, Alexander H. Liu, Leonid Karlinsky, and James R. Glass. 2024. Listen, think, and understand. In The Twelfth International Conference on Learning Representations.
- Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song Wen, Ziyu Guo, Xudong Lu, Shuai Ren, Yafei Wen, Xiaoxin Chen, Xiangyu Yue, Hongsheng Li, and Yu Qiao. 2023. ImageBind-LLM: Multi-modality instruction tuning. Preprint, arxiv:2309.03905 [cs, eess].
- Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 2023. 3d-LLM: Injecting the 3d world into large language models. Preprint, arxiv:2307.12981 [cs].
- Joy Hsu, Jiayuan Mao, and Jiajun Wu. 2023. NS3d: Neuro-symbolic grounding of 3d objects and relations. Preprint, arxiv:2303.13483 [cs].
- Haifeng Huang, Zehan Wang, Rongjie Huang, Luping Liu, Xize Cheng, Yang Zhao, Tao Jin, and Zhou Zhao. 2023a. Chat-3d v2: Bridging 3d scene and large language models with object identifiers. *Preprint*, arxiv:2312.08168 [cs].
- Haifeng Huang, Zehan Wang, Rongjie Huang, Luping Liu, Xize Cheng, Yang Zhao, Tao Jin, and Zhou Zhao. 2023b. Chat-scene: Bridging 3d scene and large language models with object identifiers. arXiv preprint arXiv:2312.08168.
- Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. 2024a. An embodied generalist agent in 3d world. Preprint, arxiv:2311.12871 [cs].

- Shijia Huang, Yilun Chen, Jiaya Jia, and Liwei Wang. 2022. Multi-view transformer for 3d visual grounding. Preprint, arxiv:2204.02174 [cs].
- Zhening Huang, Xiaoyang Wu, Xi Chen, Hengshuang Zhao, Lei Zhu, and Joan Lasenby. 2024b. Openins3d: Snap and lookup for 3d open-vocabulary instance segmentation. European Conference on Computer Vision.
- Ayush Jain, Nikolaos Gkanatsios, Ishita Mediratta, and Katerina Fragkiadaki. 2022. Bottom up top down detection transformers for language grounding in images and point clouds. Preprint, arXiv:2112.08879.
- Zhao Jin, Munawar Hayat, Yuwei Yang, Yulan Guo, and Yinjie Lei. 2023. Context-aware alignment and mutual masking for 3d-language pre-training. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10984–10994.
- Alon Lavie and Abhaya Agarwal. 2007. Meteor: an automatic metric for mt evaluation with high levels of correlation with human judgments. In Proceedings of the Second Workshop on Statistical Machine Translation.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In International Conference on Machine Learning, pages 19730-19742. PMLR.
- Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. 2023. VILA: On pre-training for visual language models. Preprint, arXiv:2312.07533.
- Dingning Liu, Xiaoshui Huang, Yuenan Hou, Zhihui Wang, Zhenfei Yin, Yongshun Gong, Peng Gao, and Wanli Ouyang. 2024a. Uni3d-llm: Unifying point cloud perception, generation and editing with large language models. arXiv preprint arXiv:2402.03327.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024b. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 26296-26306.
- Haotian Liu, Chunyuan Li, Oingyang Wu, and Yong Jae Lee. 2024c. Visual instruction tuning. Advances in Neural Information Processing Systems, 36.
- Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing Ren, Hao Shen, Huaxia Xia, and Si Liu. 2022. 3d-sps: Single-stage 3d visual grounding via referred point progressive selection. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), page 16433-16442. IEEE.
- Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang. 2023. SQA3d: Situated question answering in 3d scenes. Preprint, arxiv:2210.07474 [cs].

794

- 810 811
- 812 813
- 815

816

818 819

821

- 823
- 824 825
- 827

828

834

- 836
- 837 838

841 842

843

- 847

- Meta AI. 2023. Meta-Llama-3-8B-Instruct. https://huggingface.co/meta-llama/ Meta-Llama-3-8B-Instruct. Accessed: YYYY-MM-DD.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In ACL.
- Zekun Qi, Runpei Dong, Shaochen Zhang, Haoran Geng, Chunrui Han, Zheng Ge, Li Yi, and Kaisheng Ma. 2024. Shapellm: Universal 3d object understanding for embodied interaction. arXiv preprint arXiv:2402.17766.
- Jonas Schult, Francis Engelmann, Alexander Hermans, Or Litany, Siyu Tang, and Bastian Leibe. 2023. Mask3d: Mask transformer for 3d semantic instance segmentation. Preprint, arxiv:2210.03105 [cs].
- Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, Yuxuan Wang, and Chao Zhang. 2024. video-SALMONN: Speech-enhanced audio-visual large language models. In Forty-first International Conference on Machine Learning.
- Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. 2020. Fourier features let networks learn high frequency functions in low dimensional domains. Preprint, arXiv:2006.10739.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. 2024a. SALMONN: Towards generic hearing abilities for large language models. In The Twelfth International Conference on Learning Representations.
- Yuan Tang, Xu Han, Xianzhi Li, Qiao Yu, Yixue Hao, Long Hu, and Min Chen. 2024b. Minigpt-3d: Efficiently aligning 3d point clouds with large language models using 2d priors. arXiv preprint arXiv:2405.01413.
- Sebastian Thrun. 2008. Simultaneous Localization and Mapping, pages 13-41. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Ozan Unal, Christos Sakaridis, Suman Saha, and Luc Van Gool. 2024. Four ways to improve verbo-visual fusion for dense 3d visual grounding. Preprint, arxiv:2309.04561 [cs].
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In CVPR.
- Zehan Wang, Haifeng Huang, Yang Zhao, Linjun Li, Xize Cheng, Yichen Zhu, Aoxiong Yin, and Zhou Zhao. 2023a. 3drp-net: 3d relative positionaware network for 3d visual grounding. Preprint, arxiv:2307.13363 [cs].

Zehan Wang, Haifeng Huang, Yang Zhao, Ziang Zhang, and Zhou Zhao. 2023b. Chat-3d: Data-efficiently tuning large language model for universal dialogue of 3d scenes. Preprint, arxiv:2308.08769 [cs].

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

- Tung-Yu Wu, Sheng-Yu Huang, and Yu-Chiang Frank Wang. 2024. Data-efficient 3d visual grounding via order-aware referring. Preprint, arXiv:2403.16539.
- Yanmin Wu, Xinhua Cheng, Renrui Zhang, Zesen Cheng, and Jian Zhang. 2023. Eda: Explicit text-decoupling and dense alignment for 3d visual grounding. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), page 19231-19242. IEEE.
- Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. 2023. Pointllm: Empowering large language models to understand point clouds. arXiv preprint arXiv:2308.16911.
- Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. 2023. Ulip: Learning a unified representation of language, images, and point clouds for 3d understanding. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 1179–1189.
- Jianing Yang, Xuweiyi Chen, Shengyi Qian, Nikhil Madaan, Madhavan Iyengar, David F Fouhey, and Joyce Chai. 2024. Llm-grounder: Open-vocabulary 3d visual grounding with large language model as an agent. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 7694–7701. IEEE.
- Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. 2022. Point-bert: Pretraining 3d point cloud transformers with masked point modeling. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 19313-19322.
- Hang Zhang, Xin Li, and Lidong Bing. 2023a. Video-LLaMA: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858.
- Yiming Zhang, ZeMing Gong, and Angel X. Chang. 2023b. Multi3drefer: Grounding text description to multiple 3d objects. Preprint, arxiv:2309.05251 [cs].
- Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang. 2023. Uni3d: Exploring unified 3d representation at scale. arXiv preprint arXiv:2310.06773.
- Zhu Ziyu, Ma Xiaojian, Chen Yixin, Deng Zhidong, Huang Siyuan, and Li Qing. 2023. 3d-vista: Pretrained transformer for 3d vision and text alignment. In ICCV.

- 902
- 903
- 904 905
- 906
- 908
- 910
- 911
- 912 913
- 914
- 915
- 916
- 917

942

918 919 921

Annotation Procedure for Α **DetailedScanRefer**

A.1 Comparison between rendered image and real-world photo

We explored the use of rendered images for annotation but found the image quality lacking compared to real-world photos retrieved from the Scannet dataset using the Camera Pose Matching algorithm. The rendered images often suffer from poor lighting, texture quality, and geometric accuracy, making them less suitable for precise annotations. By comparing these renderings with ScanNet photos, the limitations of synthetic data are clear, highlighting the need for higher-fidelity imagery or realworld data for accurate annotation tasks.

A.2 example of pixel-level visibility masks

A.3 GPT-40 Annotation Prompt

This is the prompt we used for GPT annotation:

"You are a helpful assistant designed to output JSON. The task is to identify all mentioned objects in the image and add the matching obj id to the given description. The OBJID is shown in red font, and it should be annotated at the cen-923 tre of the object. Remember, please return both the <input_description> and the <augmented_description> with obj 926 id added. You should not modify the <input_description>. Only add the <OB-928 929 JID> after the object entity if you can recognize both the object and the red an-930 notation clearly in the image. Also, if 931 you cannot recognize ALL of the objects AND ALL of their corresponding red an-933 notation in the description, simply output "NAN" in the "augmented description". 935 An example is here: "input_description: This is a brown chair. it is at a high ta-937 ble. augmented_description: This is a 938 brown chair <OBJ003>. it is at a high table <OBJ012>."

> This is the prompt format we used for the DetailedScanrefer Dataset:

"According to the provided description, 943 <input_description>, please append the correct object ID after each object mentioned in the description."

The <input description> refers to the original 947 ScanRefer description. The annotations generated 948 by GPT-40 serve as the reference captions for each 949 corresponding question. 950

951

952

953

954

955

956

957

A.4 Dataset Statistics

We provide the dataset statistics in Tab.5. The numbers of descriptions before and after each processing step are shown.

Descript	Description						
Before pr	Before processing						
Inconsist	ent first O	bjId	13,836				
NaN valu	ies		2,191				
No ObjId	l		154				
Invalid O	bjId range	e	6				
After pro	16,151						
Table 5: I	Data Proce	ssing S	Statistics				
Rating	Count	Perc	entage				
1	28	0.	09%				
2	3291	10	.18%				
3	7092	21	.93%				
4	17	.38%					
5	5 16310 5						
Total	32338	1	00%				

Table 6: Distribution of GPT-40 Ratings

A.5 Dataset Quality Evaluation

We also evaluate the annotation quality automatically using GPT40 with the following prompt:

"You are tasked with evaluating the ac-958 curacy and completeness of text anno-959 tations provided for objects in an image. 960 Some objects in the image is labeled with 961 an object ID (e.g., <OBJ014>), and these 962 IDs are referenced in the text annotations. 963 Your goal is to ensure that every object 964 mentioned in the text annotation has a ac-965 curate corresponding red-text annotation 966 in the image. First, verify that all ob-967 jects mentioned in the text are annotated 968 in the image. Second, ensure that the 969 object descriptions in the text correctly 970 match the labeled objects in the image 971 in terms of type, appearance, and loca-972 tion. After reviewing, provide a rating 973 between 1 and 5, where 1 represents poor 974



Figure 5: Top row: Scannet Photos. Bottom row: Corresponding Rendered Images. Each pair of images corresponds to the same ScanRefer description.

annotation quality and 5 represents excellent quality. The rating should consider whether all objects mentioned in the text are annotated in the image, and whether the descriptions are accurate."

975

976

979

982

983

985

991

992

993

995

997

1001

1002

1003

1004

1005

1006

1008

We submitted both the annotated photo and its detailed description to GPT-40, requesting it to rate their consistency and accuracy on a scale from 1 to 5. An average rating of 3.31 was achieved across 32,338 descriptions. Specifically, the distribution of ratings is presented in Tab.6. Two examples are shown in Fig.7, where the annotation, rating, and its reasoning are illustrated.

B Additional Results on Scan2Cap

In this appendix, we provide a detailed explanation of the modifications made to the Scan2Cap dataset for our viewpoint-aware captioning task. We also discuss the implications for evaluation and how these changes affect comparability with existing models.

Scan2Cap (Chen et al., 2021) is a captioning dataset generated based on **ScanRefer** (Chen et al., 2020), which provides natural language descriptions of objects within 3D indoor scenes from the ScanNet dataset (Dai et al., 2017).

Each description in ScanRefer is associated with a specific camera pose. By reusing these camera poses, we reconstruct viewpoints for the Scan2Cap dataset, making it **viewpoint-aware**. This approach enhances the dataset by incorporating spatial context and specific viewpoints, providing a more comprehensive captioning task.

The Scan2Cap dataset utilizes a set of predefined prompts to guide the captioning task. Notably,

these prompts are used in the original Chat-3D-
v2 task (Huang et al., 2023a). The prompts are
designed to elicit detailed descriptions of objects
and their spatial relationships within a scene. The
prompts include:1009
10101010
10111012
1012

1. "Begin by detailing the visual aspects of the 1014 <id> before delving into its spatial context 1015 among other elements within the scene." 1016 2. "First, depict the physical characteristics of 1017 the <id>, followed by its placement and inter-1018 actions within the surrounding environment." 1019 3. "Describe the appearance of the <id>, then 1020 elaborate on its positioning relative to other 1021 objects in the scene." 1022 4. "Paint a picture of the visual attributes of <id>, 1023 then explore how it relates spatially to other 1024 elements in the scene." 5. "Start by articulating the outward features of 1026 the <id>, then transition into its spatial align-1027 ment within the broader scene." 1028 6. "Provide a detailed description of the appear-1029 ance of <id> before analyzing its spatial connections with other elements in the scene." 1031 7. "Capture the essence of the appearance of 1032 <id>, then analyze its spatial relationships 1033 within the scene's context." 1034 8. "Detail the physical characteristics of the <id> 1035 and subsequently examine its spatial dynam-1036

ics amidst other objects in the scene."

1037



Figure 6: Pixel-level segmentation masks for object visibility checks. Top: original photos. Middle: corresponding segmentation masks. Bottom: annotated photos with object identifiers.

9. "Describe the visual traits of <id> first, then elucidate its spatial arrangements in relation to neighboring elements."

1038

1039

1040

1041

1042

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1061

10. "Begin by outlining the appearance of <id>, then proceed to illustrate its spatial orientation within the scene alongside other objects."

An example entry from the original Scan2Cap dataset is provided below. Multiple reference captions correspond to a single prompt, offering varied descriptions of the object.

Prompt: First, depict the physical characteristics of the <OBJ014>, followed by its placement and interactions within the surrounding environment.

The corresponding reference captions are as follows:

- There are brown wooden cabinets. Placed on the side of the kitchen.
- There is a set of bottom kitchen cabinets in the room. It has a microwave in the middle of it.
- There is a set of bottom kitchen cabinets in the room. There is a microwave in the middle of them.

Brown kitchen cabinets, the top is decorated with marble layers, and it is placed on the left in the direction of view. On the right, there are four brown chairs.

1066

1067

1069

1070

1071

1072

1074

1075

1076

1081

1082

• The kitchen cabinets are located along the right wall. They are below the countertop. The kitchen cabinets are located to the right of the table and chairs.

In our modified viewpoint-aware Scan2Cap dataset, each question is associated with a specific viewpoint (camera pose). Each viewpoint corresponds to one correct reference caption. The dataset includes camera parameters (position, rotation, lookat) for each entry. Example entries:

Prompt: "First, depict the physical char-
acteristics of the <OBJ014>, followed by
its placement and interactions within the
surrounding environment.",1077
1078

Given a particular camera pose it would have one reference caption corresponding to the origin Scanrefer Dataset:

"There is a set of bottom kitchen cabinets in the room. There is a microwave in the middle of them."



GPT Annotation: "this is a wicker cabinet <OBJ019>. it is to the right of the boxes <OBJ039>."

Rating: 5/5

Reasoning: The annotations are correct. The wicker cabinet labeled as <OBJ019> matches the object in the image. The boxes labeled <OBJ039> are also correctly identified and are to the left of the wicker cabinet, which aligns with the described spatial relationship.



GPT Annotation: "this is a brown cabinet <OBJ027>. it is to the right of a picture <OBJ032>"

Rating: 1/5

Reasoning: The object IDs are incorrectly attributed in the description. The brown cabinet is identified as <OBJ027>, which matches the visible cabinet in the image. However, the picture is labeled as <OBJ032> in the image, but spatially it does not appear to be to the left of the cabinet. Additionally, the context does not logically support the described spatial relationship.

Figure 7: Comparison of two rating scenarios in different scenes.

The modification of the dataset has significant implications for evaluation: The evaluation method changes due to the dataset modification. Each question under a specific viewpoint has a single correct reference caption. This differs from the original method, which averaged over multiple reference captions. Consequently, the modified dataset and evaluation method are **not directly comparable** to models trained on the original Scan2Cap dataset.

C Detailed Comparison of GeVLM and other models

		@0.25	IoU	@0.5 IoU					
System	CIDER	BLEU-4	METEOR	CIDER	BLEU-4	METEOR			
Chat-3D-v2* (Huang et al., 2023a)	72.20	11.28	18.89	68.63	10.46	18.23			
Ours	67.94	11.19	19.06	64.47	10.44	18.38			

Table 7: Evaluation results on Scan2Cap validation set at IoU thresholds 0.25 and 0.50.

Table 8: F1 scores at 0.5 IoU on Multi3DRefer (Zhang et al., 2023b) validation set. ZT, ST, and MT refer to zero, single, and multiple target objects in the scene referenced by each description. "D" refers to distracting objects of the same semantic class.

	ZT w/o D	ZT w/ D	ST w/o D	ST w/ D	MT	All
Chat-3D v2* (Huang et al., 2023a)	90.7	62.2	64.3	33.0	42.1	44.9
Chat-Scene	90.3	62.6	75.9	44.5	41.1	52.4
GeVLM (Chat-3D v2) (Ours)	90.3	68.0	68.0	36.5	36.7	46.1
GeVLM (Chat-Scene) (Ours)	92.0	63.2	75.7	46.5	36.8	52.4

Table 9: Results on language tasks using ScanQA validation set and SQA3D test set. B1 to B4 represents BLEU-1 to 4, M for METEOR, C for CIDEr, and R for ROUGE-L. The Chat-3D v2 is the reproduced results which is slightly better than the reported numbers in the original paper.

System				Scar	ıQA				SQA3D						
System	B1	B2	B3	B4	М	С	R	EM	What	Is	How	Can	WhichOthersAvg		sAvg
LL3DA (Chen et al., 2024a)	_	-	-	13.5	15.9	76.8	37.3	-	_	-	_	_	_	_	_
LEO (Huang et al., 2024a)	-	_	_	13.2	20.0	101.4	4 4 9.2	24.5	_	_	-	-	-	-	50.0
Scene-LLM (Fu et al., 2024)	42.2	26.4	18.7	11.7	15.8	80.0	35.9	25.6	40.0	69.2	42.8	70.8	46.6	52.5	53.6
3D-LLM (BLIP2-flant5) (Hong et al., 2023)	39.3	25.2	18.4	12.0	14.5	69.4	35.7	20.5	_	_	_	-	-	-	-
Chat-3D v2* (Huang et al., 2023a)	42.3	28.1	19.6	13.4	18.0	88.9	42.1	22.4	43.9	66.0	52.5	66.3	46.4	50.2	52.5
Chat-scene	43.2	29.1	20.6	14.3	18.0	87.7	41.6	21.6	45.4	67.0	52.0	69.5	49.9	55.0	54.6
GeVLM (Chat-3D v2) (Ours)	42.4	28.7	21.3	15.4	18.1	90.5	41.8	21.7	44.1	68.6	52.3	62.7	45.6	55.8	53.5
GeVLM (Chat-scene) (Ours)	44.7	29.5	20.8	14.3	18.3	89.7	42.7	21.9	48.2	71.0	52.7	68.9	47.6	57.8	56.5

Catagory	System	Uni	que	Mult	iple	Ove	rall
Category	System	Acc@0.25	Acc@0.5	Acc@0.25	Acc@0.5	Acc@0.25	Acc@0.5
	ScanRefer (Chen et al., 2020)	76.33	53.51	32.73	21.11	41.19	27.40
	MVT (Huang et al., 2022)	77.67	66.45	31.92	23.30	39.43	33.26
	3D-SPS (Luo et al., 2022)	84.12	66.72	40.32	29.82	48.36	36.98
	ViL3DRel (Chen et al., 2022c)	81.58	68.62	40.30	30.71	47.94	37.73
	BUTD-DETR (Jain et al., 2022)	84.20	66.30	46.60	35.10	52.20	39.80
	HAM (Chen et al., 2022b)	79.24	67.86	41.46	34.03	48.79	40.60
Free and Madal	3DRP-Net (Wang et al., 2023a)	83.13	67.74	42.14	31.95	50.10	38.90
Expert Model	EDA (Wu et al., 2023)	85.76	68.57	49.13	37.64	54.59	42.26
	M3DRef-CLIP (Zhang et al., 2023b)	85.30	77.20	43.80	36.80	51.90	44.70
	ConcreteNet (Unal et al., 2024)	86.40	82.05	42.41	38.39	50.61	46.53
	DORa (Wu et al., 2024)	-	-	-	-	52.80	44.80
	3D-VLP (Jin et al., 2023)	84.23	64.61	43.51	33.41	51.41	39.46
	3D-VisTA (Ziyu et al., 2023)	81.60	75.10	43.70	39.10	50.60	45.80
	3DJCG (Cai et al., 2022)	83.47	64.34	41.39	30.82	49.56	37.33
3D Grounding + 3D Captioning	D3Net (Chen et al., 2022a)	-	72.04	-	30.05	-	37.87
3D Grounding + 3D Captioning +3D Q&A	GeVLM (Ours)	82.00	75.70	39.00	34.70	46.90	42.30

Table 10: Performance comparison of models on Scanrefer.

D Qualitative result of Mask3D

1098

and validation splits, including counts for $IoU \ge$ 1101 0.25 and $IoU \ge 0.50$, along with maximum IoU 1102 rates, to further demonstrate the quality and com- 1103

1099The table summarizing key statistics and metrics1100for Mask3D's performance across both the training

prehensiveness of the generated proposals.

1105	In the training process, we only use 32,338 anno-
1106	tations that meet the strict criterion of $IoU \ge 0.75$
1107	with ground truth objects, ensuring that only highly
1108	accurate object proposals are retained. This high
1109	threshold reflects the precision and relevance of the
1110	dataset for effective downstream tasks.

Metric	Train Split Count	Validation Split Count
Total Count (Original ScanRefer Dataset)	36,665	9,508
$IoU \ge 0.25$ Count	36,187	8,924
$IoU \ge 0.50$ Count	35,061	8,168
Max IoU@0.25	0.9870	0.9386
Max IoU@0.50	0.9563	0.8591

Table 11: Summary of Mask3D-generated dataset metrics for training and validation splits.

E Qualitative examples

1112Qualitative examples showcasing the advantage of1113GeVLM are given in Figure 8, where four types1114of confusions caused by viewpoint consistency are1115provided.



(A). <u>Description</u>: This is a black computer monitor. The black computer monitor sits on the far right of the desk.



(B). <u>Description</u>: The chair is blue with a white shirt thrown over the back. it is the chair on the right that is closes to the far wall



(C). *Description:* There is a light grey pillow on the bed. it is smaller than the other pillows and in front of the red pillow.



(D). <u>Description</u>: The clothes dryer is in the northeast corner of the room. the clothes dryer has a white color and a half circle mirror in the center

Figure 8: Comparison between GeVLM and baseline (Chat-3D v2) on viewpoint-related examples with potential ambiguities, including: (A) left/right, (B) near/far, (C) front/back and (D) north/south/east/west. The description is associated with a specific viewpoint and hence becomes confusing in other viewpoints.

Discussion of License F

1116

1121

GPT-40 was used to generate annotations for im-1117 ages. According to the terms of use, the output is 1118 owned by the users and can be used for academic 1119 purposes for DetailedScanRefer as long as they are 1120 not "to develop any artificial intelligence models that compete with OpenAI products and services". 1122 We believe that the generated content is allowed 1123 to be distributed. We will also clearly state in our 1124 licence upon public release of the dataset to refer-1125 ence the OpenAI licence. 1126